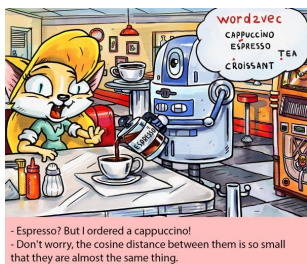**Going from Text to Knowledge Graphs:**
**Putting Natural Language Processing and Graph Databases to Work**

Dr. Clair J. Sullivan
Graph Data Science Advocate
clair.sullivan@neo4j.com
@CJLovesData1

---



2   https://www.kdnuggets.com/2017/04/cartoon-word2vec-espresso-cappuccino.html

---

**Outline**

- Introduction to knowledge graphs
- Introduction to NLP (as applied to knowledge graphs)
- How to use NLP to create a knowledge graph
- Software setup
- DEMO

3

**All materials (including the slides!) for this workshop are available on the workshop GitHub repo:**

https://github.com/cj2001/odsc_east_kg_2021

4

**Speaking of which, let's get ready!**

git clone https://github.com/cj2001/odsc_east_kg_2021

docker-compose build

5

# Two Key Concepts

6

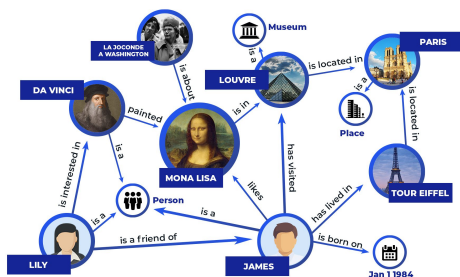**1. There is no proverbial "silver bullet" with NLP**

**2. The quality of what you get out of a knowledge graph depends on the quality of what you put into it**
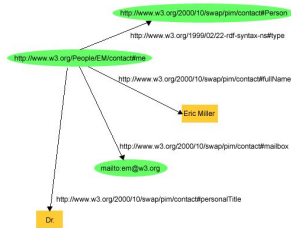
# Introduction to knowledge graphs

- "Things not strings"
- What knowledge graphs are useful for
  - Search
  - Question answering
  - Recommendation engine
- Can be generated a lot of different ways
  - Resource Description Framework (RDF)
  - Co-occurrence
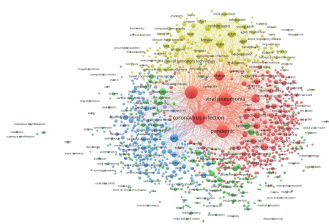  - Subject-Verb-Object (SVO)

10

---

# RDF triples



http://www.w3.org/2000/10/swap/pim/contact#Person

http://www.w3.org/1999/02/22-rdf-syntax-ns#type

http://www.w3.org/People/EM/contact#me

http://www.w3.org/2000/10/swap/pim/contact#fullName

Eric Miller

http://www.w3.org/2000/10/swap/pim/contact#mailbox

mailto:em@w3.org

http://www.w3.org/2000/10/swap/pim/contact#personalTitle

Dr.

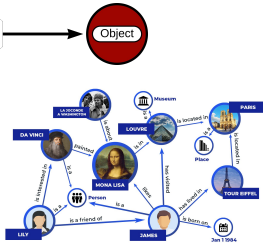11   https://en.wikipedia.org/wiki/Resource_Description_Framework#Examples

---

# Word co-occurrence



12   https://covid19bibio.com/2020/04/28/keyword-co-occurrence-network-graph-for-the-overall-research-field-on-covid-19-up-to-april-27th-2020/

## SVO triples



---

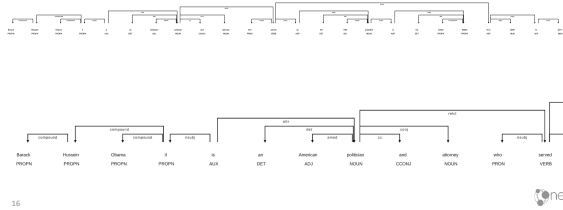## NLP considerations for knowledge graph creation

- Named Entity Recognition (NER)
- SVO / SPO triples
  - …but verbs can be difficult, as you will see!
- Very language dependent
- Very topic-area dependent

---

Barack Hussein Obama II is an American politician and attorney who served as the 44th president of the United States from 2009 to 2017.



16

---

Barack Hussein Obama II is an American politician and attorney who served as the 44th president of the United States from 2009 to 2017.

| Text | Lemma | Tag | POS | DEP | is_stop |
|---|---|---|---|---|---|
| Barack | Barack | NNP | PROPN | compound | FALSE |
| Hussein | Hussein | NNP | PROPN | compound | FALSE |
| Obama | Obama | NNP | PROPN | compound | FALSE |
| II | II | NNP | PROPN | nsubj | FALSE |
| is | be | VBZ | AUX | ROOT | TRUE |
| an | an | DT | DET | det | TRUE |
| American | american | JJ | ADJ | amod | FALSE |
| politician | politician | NN | NOUN | attr | FALSE |
| and | and | CC | CCONJ | cc | TRUE |
| attorney | attorney | NN | NOUN | conj | FALSE |
| who | who | WP | PRON | nsubj | TRUE |
| served | serve | VBD | VERB | relcl | FALSE |
| as | as | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| 44th | 44th | JJ | ADJ | amod | FALSE |
| president | president | NN | NOUN | pobj | FALSE |
| of | of | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| United | United | NNP | PROPN | compound | FALSE |
| States | States | NNP | PROPN | pobj | FALSE |
| from | from | IN | ADP | prep | TRUE |
| 2009 | 2009 | CD | NUM | pobj | FALSE |
| to | to | IN | ADP | prep | TRUE |
| 2017 | 2017 | CD | NUM | pobj | FALSE |
| . | . | . | PUNCT | punct | FALSE |

17

---

Barack Hussein Obama II is an American politician and attorney who served as the 44th president of the United States from 2009 to 2017.

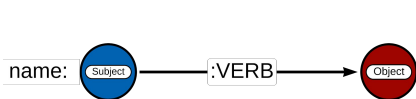| Text | Lemma | Tag | POS | DEP | is_stop |
|---|---|---|---|---|---|
| Barack | Barack | NNP | PROPN | compound | FALSE |
| Hussein | Hussein | NNP | PROPN | compound | FALSE |
| Obama | Obama | NNP | PROPN | compound | FALSE |
| II | II | NNP | PROPN | nsubj | FALSE |
| is | be | VBZ | AUX | ROOT | TRUE |
| an | an | DT | DET | det | TRUE |
| American | american | JJ | ADJ | amod | FALSE |
| politician | politician | NN | NOUN | attr | FALSE |
| and | and | CC | CCONJ | cc | TRUE |
| attorney | attorney | NN | NOUN | conj | FALSE |
| who | who | WP | PRON | nsubj | TRUE |
| served | serve | VBD | VERB | relcl | FALSE |
| as | as | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| 44th | 44th | JJ | ADJ | amod | FALSE |
| president | president | NN | NOUN | pobj | FALSE |
| of | of | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| United | United | NNP | PROPN | compound | FALSE |
| States | States | NNP | PROPN | pobj | FALSE |
| from | from | IN | ADP | prep | TRUE |
| 2009 | 2009 | CD | NUM | pobj | FALSE |
| to | to | IN | ADP | prep | TRUE |
| 2017 | 2017 | CD | NUM | pobj | FALSE |
| . | . | . | PUNCT | punct | FALSE |

18

Barack Hussein Obama II is an American politician and attorney who served as the 44th president of the United States from 2009 to 2017.

| Text | Lemma | Tag | POS | DEP | is_stop |
|---|---|---|---|---|---|
| Barack | Barack | NNP | PROPN | compound | FALSE |
| Hussein | Hussein | NNP | PROPN | compound | FALSE |
| Obama | Obama | NNP | PROPN | compound | FALSE |
| II | II | NNP | PROPN | nsubj | FALSE |
| is | be | VBZ | AUX | ROOT | TRUE |
| an | an | DT | DET | det | TRUE |
| American | american | JJ | ADJ | amod | FALSE |
| politician | politician | NN | NOUN | attr | FALSE |
| and | and | CC | CCONJ | cc | TRUE |
| attorney | attorney | NN | NOUN | conj | FALSE |
| who | who | WP | PRON | nsubj | TRUE |
| served | serve | VBD | VERB | relcl | FALSE |
| as | as | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| 44th | 44th | JJ | ADJ | amod | FALSE |
| president | president | NN | NOUN | pobj | FALSE |
| of | of | IN | ADP | prep | TRUE |
| the | the | DT | DET | det | TRUE |
| United | United | NNP | PROPN | compound | FALSE |
| States | States | NNP | PROPN | pobj | FALSE |
| from | from | IN | ADP | prep | TRUE |
| | 2009 | 2009 | CD | NUM | pobj | FALSE |
| to | to | IN | ADP | prep | TRUE |
| | 2017 | 2017 | CD | NUM | pobj | FALSE |
| . | . | | PUNCT | punct | FALSE |

---

## Detailed knowledge graph data model

name: (Subject) :VERB → (Object)

name:
node_labels (*):
description:
url:
word_vec:

---

## An introduction to the tools we will use today

- spacy
- Wikipedia Python package
- Google Knowledge Graph
- Neo4j
    - Awesome Procedures on Cypher (APOC)
    - Graph Data Science (GDS) Library
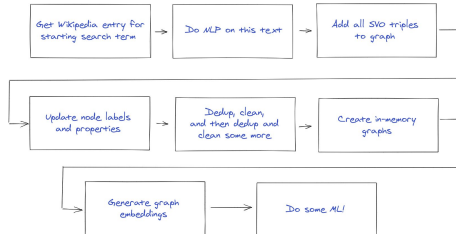    - Cypher
        - No Cypher knowledge is assumed!

```
{...
  "@type": "ItemList",
  "itemListElement": [
    {
      "@type": "EntitySearchResult",
      "result": {
        "@id": "kg:/m/0dl567",
        "name": "Taylor Swift",
        "@type": [
          "Thing",
          "Person"
        ],
...
        "detailedDescription": {
          "articleBody": "Taylor Alison Swift is an American singer-songwriter and
actress. Raised in Wyomissing, Pennsylvania, she moved to Nashville, Tennessee, at the
age of 14 to pursue a career in country music. ",
          "url": "http://en.wikipedia.org/wiki/Taylor_Swift",
...
}
```
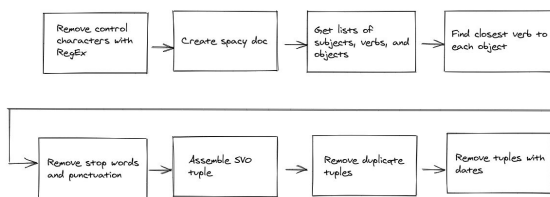
22

neo4j

---

## Overview of workflow

```
Get Wikipedia entry for      →   Do NLP on this text   →   Add all SVO triples
starting search term                                        to graph

Update node labels           →   Dedup, clean,          →   Create in-memory
and properties                   and then dedup and          graphs
                                 clean some more

Generate graph               →   Do some ML!
embeddings
```

23

neo4j

---

## NLP workflow

```
Remove control        →   Create spacy doc   →   Get lists of        →   Find closest verb to
characters with                                   subjects, verbs, and     each object
RegEx                                             objects

Remove stop words     →   Assemble SVO       →   Remove duplicate    →   Remove tuples with
and punctuation           tuple                   tuples                   dates
```

24

neo4j

## Detailed knowledge graph data model

name: [Subject] —— :VERB ——▶ [Object]

name:
node_labels (*):
description:
url:
word_vec:

## Software setup

- You will need access to Docker, CLI, and a browser
- Clone the GitHub repository:
  - https://github.com/cj2001/odsc_east_kg_2021
- Google KG API key
  - https://developers.google.com/knowledge-graph/prereqs
  - Get this key and then save it in /notebooks/.api_key

Home > Search Central > Knowledge Graph Search API                    Rate and review  👍 👎

## Prerequisites

Before you can start coding your first client application, there are a few things you need to do, if you haven't done them already.

### Get a Google Account

You need a Google Account in order to create a project in the Google API Console. If you already have an account, then you're all set.

You may also want a separate Google Account for testing purposes.

### Create a project for your client

Before you can send requests to Google Knowledge Graph Search API, you need to tell Google about your client and activate access to the API. You do this by using the Google API Console to create a project, which is a named collection of settings and API access information, and register your application.

To get started using Google Knowledge Graph Search API, you need to first use the setup tool, which guides you through creating a project in the Google API Console, enabling the API, and creating credentials.

If you haven't done so already, create your application's API key by clicking **Create credentials > API key**. Next, look for your API key in the **API keys** section.

## Software setup

- At this point you should have:
  - Google Knowledge Graph API key
  - Repo cloned
  - Docker container built

---

**Let's go!**

docker-compose up

---

# 1. There is no proverbial "silver bullet" with NLP

## 2. The quality of what you get out of a knowledge graph depends on the quality of what you put into it

neo4j

31

---

## Review

- Learned why knowledge graphs are cool!
- Discussed approaches to create a knowledge graph, focusing on NLP
- Created a knowledge graph using a variety of bits of information on the web
- Explored ways to make the graph more informative
  - Entity resolution / disambiguation
- Created various graph embeddings for downstream ML work

neo4j

32

---

## Key observations

- The quality of the results depends strongly on the upstream NLP
  - Should be customized to the graph
  - SME involvement helps a lot
- Think carefully about the graph model prior to creating the embeddings
  - Big implications on feature engineering

neo4j

33

## Things to try next

- Hyperparameters!!!
  - NLP
  - Graph embeddings
  - ML
- Use different properties to create graph embeddings
- Try different graph embedding methods
- Class imbalance
- Make graph bigger, include more diverse data sources

34

---

**All materials (including the slides!) for this workshop are available on the workshop GitHub repo:**

https://github.com/cj2001/odsc_east_kg_2021

35

---

# Thank you!

**@CJLovesData1**

36